

# TabularBench: Benchmarking Adversarial Robustness for Tabular Deep IIII.III

Learning in Real-world Use-cases <sup>1</sup>University of Luxembourg, <sup>2</sup>Luxembourg institute of Science and Technologies, <sup>3</sup>Riken AIP Thibault Simonetto<sup>1</sup>, Salah Ghamizi<sup>2,3</sup> and Maxime Cordy<sup>1</sup> Contact: thibault.simonetto@uni.lu

Survival

Selection

Crossover

**Model architectures** 

TabTransformer

- RLN

- STG

TabNet

- VIME





# Problem ML System Attacker objective immutable or computed $\checkmark H(x) \neq H(x + \delta)$ internally $\checkmark L_p(x, x + \delta) < \epsilon$ $\checkmark x + \delta \in \mathcal{X}_{\Omega}$ Adversarial input

2. Feature relationships

constraints are used to

reject invalid input

# $H: \mathcal{X} \to \mathcal{Y}$ the classification model $L_n$ the distance according to a p-norm $\mathcal{X}_{\Omega} \subseteq \mathcal{X}$ the subspace where x satisfies all constraints $\omega \in \Omega$

### \_\_\_\_\_\_\_ **Example of constraint:**

 $H(X) \in [0,1]^N$ 

 $int\_rate \times (1 + int\_rate)^{term}$  $installment = loan\_amount \times f$  $(1 + int_rate)^{term} - 1$ 

Respected

ML Classifier

### **Constraints grammar:**

$$\omega \coloneqq \omega_1 \wedge \omega_2 \mid \omega_1 \vee \omega_2 \mid \psi_1 \succeq \psi_2$$

$$\psi \coloneqq c \mid f_i \mid \psi_1 \oplus \psi_2 \mid x_i$$

### Known attacks produce invalid examples:

| PGD        | X |
|------------|---|
| AutoAttack | X |

|             | Max trans. | Avg trans. | Acc. creation | Age |
|-------------|------------|------------|---------------|-----|
| Adversarial | \$2500     | \$3000     | 1 year        | 22  |

Adversarial violates Avg transaction  $\leq Max$  transaction

### Constraints as a penalty function:

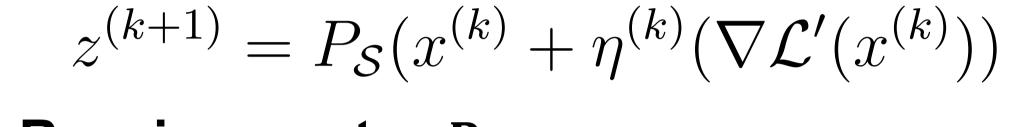
| Constraint                 | Penalty function                 | Constraint example  | Penalty function example   |
|----------------------------|----------------------------------|---|--|
| $\psi_1 = \psi_2$          | $\mid \psi_1 - \psi_2 \mid$      | $rec_per_month = record/month$  | rec_per_month - (record/month)   |
| $\psi_1 \le \psi_2$        | $max(0,\psi_1-\psi_2)$           | $open\_acc \leq total\_acc$   | $max(0, open\_acc - total\_acc)$   |
| $\psi_1 < \psi_2$          | $max(0, \psi_1 - \psi_2 + \tau)$ | open_acc < total_acc  | $max(0, open\_acc - total\_acc + 10^{-5})$   |
| $\omega_1 \wedge \omega_2$ | $\omega_1 + \omega_2$            | $((\text{term} = 36) \lor (\text{term} = 60)) \land (\text{open\_acc} \le \text{total\_acc})$ | $\min( \text{term} - 36 ,  \text{term} - 60 ) + \max(0, \text{open\_acc} - \text{total\_acc})$ |
| $\omega_1 \vee \omega_2$   | $\min(\omega_1,\omega_2)$        | $(\text{term} = 36) \lor (\text{term} = 60)$  | $\min( \text{term} - 36 ,  \text{term} - 60 )$   |



# Projected Gradient Descent (PGD) $x^{(k+1)} = P_{\mathcal{S}}(x^{(k)} + \eta \nabla l(h(x), y))$

# Gradient loss + Constraints regularization $\mathcal{L}'(x) = l(h(x), y) - \sum_{\omega_i \in \Omega} penalty(x, \omega_i)$

# **Constrained Gradient Descent**

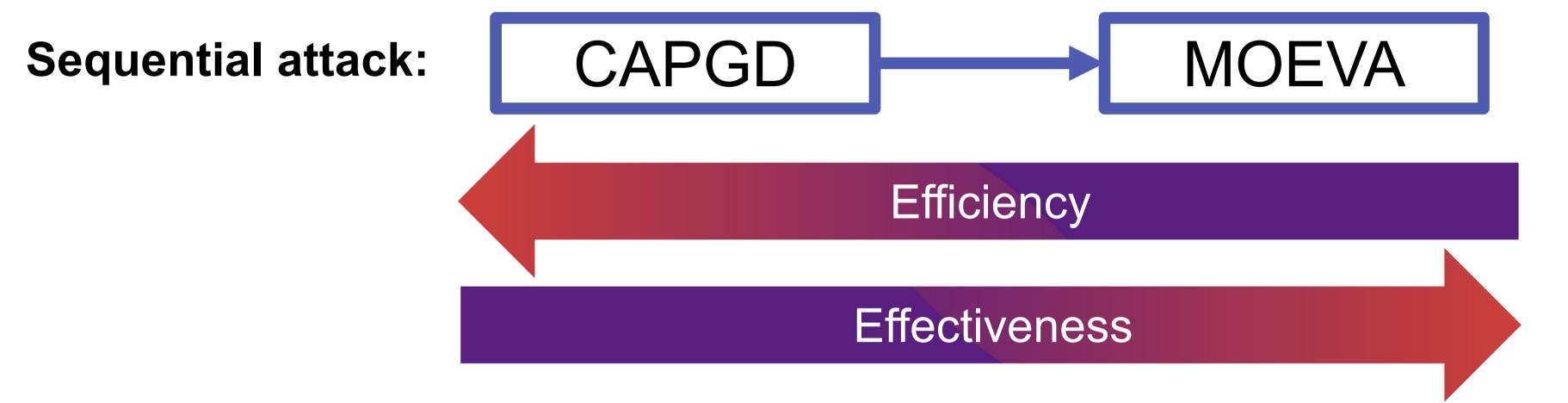


# Repair operator $R_{\Omega}$

CAPGD (white-box)



# Constrained Adaptive Attack (CAA)



# MOEVA (gray-box)

# Multi-objective genetic algorithm (NSGA-III)

 $minimise\ g_1(x) \equiv h(x)$ 

minimise  $g_2(x) \equiv L_p(x - x_0)$ 

minimise  $g_3(x) \equiv$ 

# Metrics



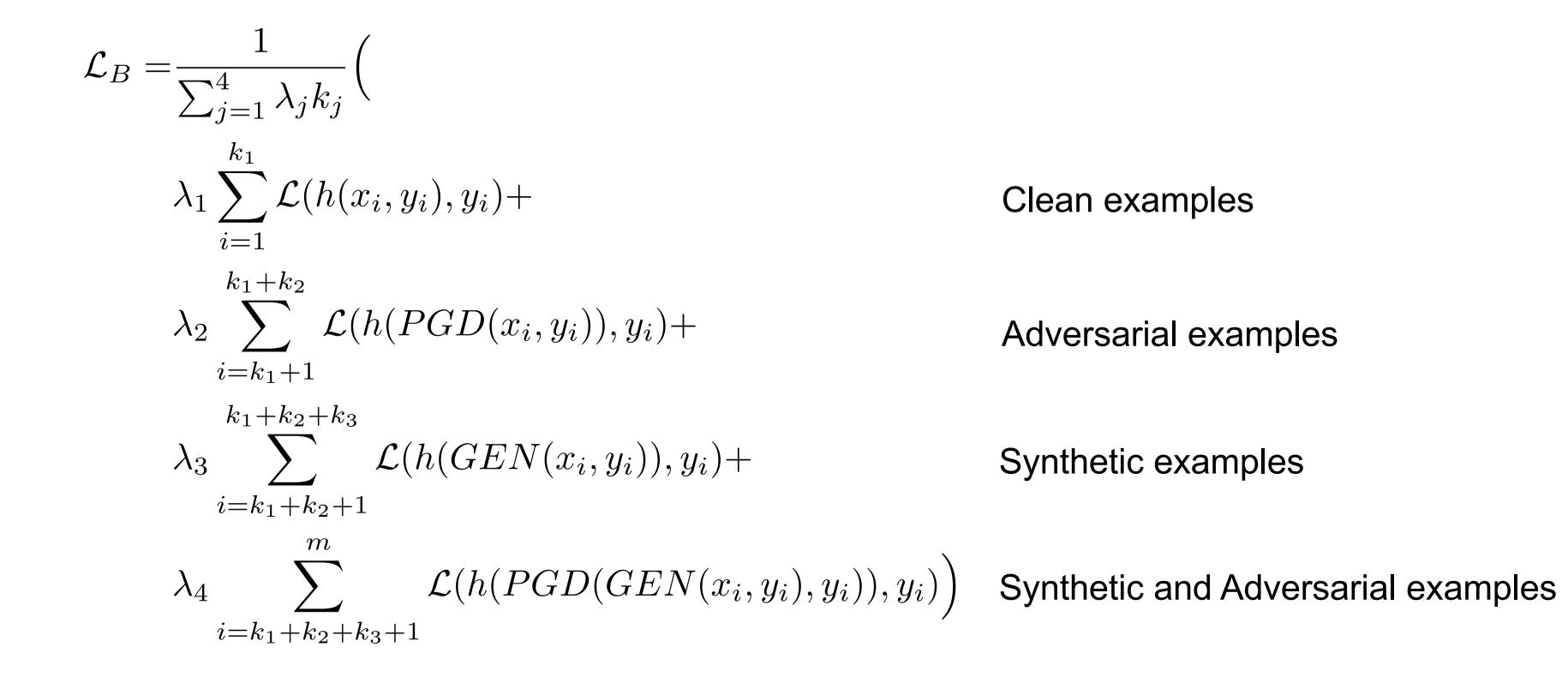
# Correctly classified adversarial examples  $Robust \ accuracy = -$ # Clean examples

# TABULARBENCH

### Regularized adversarial training

 $\hat{\mathcal{L}}_B = \frac{1}{(m-k) + \lambda k} \left( \lambda \sum_{i=1}^n \mathcal{L}(h(PGD(x_i, y_i)), y_i) + \sum_{i=k+1}^m \mathcal{L}(h(x_i), y_i) \right)$ 

### Regularized adversarial training with data augmentation



### **Training**

- Standard
- Adversarial training
- CT-GAN
  - GOGGLE

- None

TableGAN

Data augmentation

- TVAE
- WGAN CutMix

### Three domains, five datasets

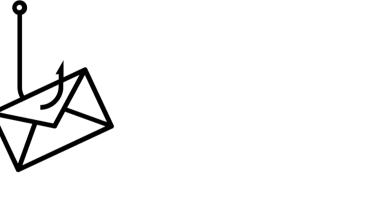


Credit scoring

Lending Club Loan Data





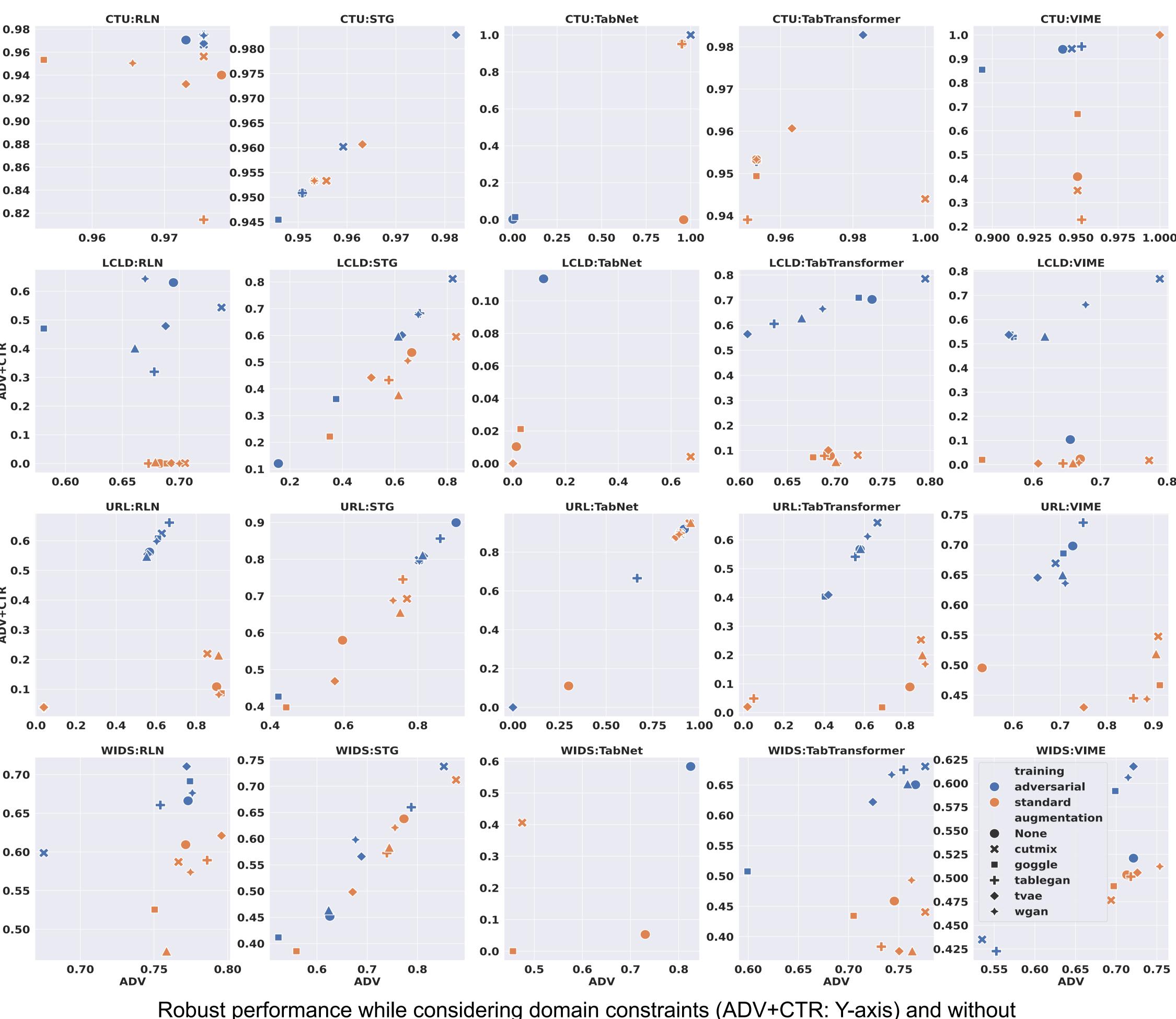




**MALWARE** 



# Results



Robust performance while considering domain constraints (ADV+CTR: Y-axis) and without (ADV: X-axis) on all our use cases confirms the relevance of studying constrained-aware attacks.

### API

### Available on PyPi, Conda and Docker

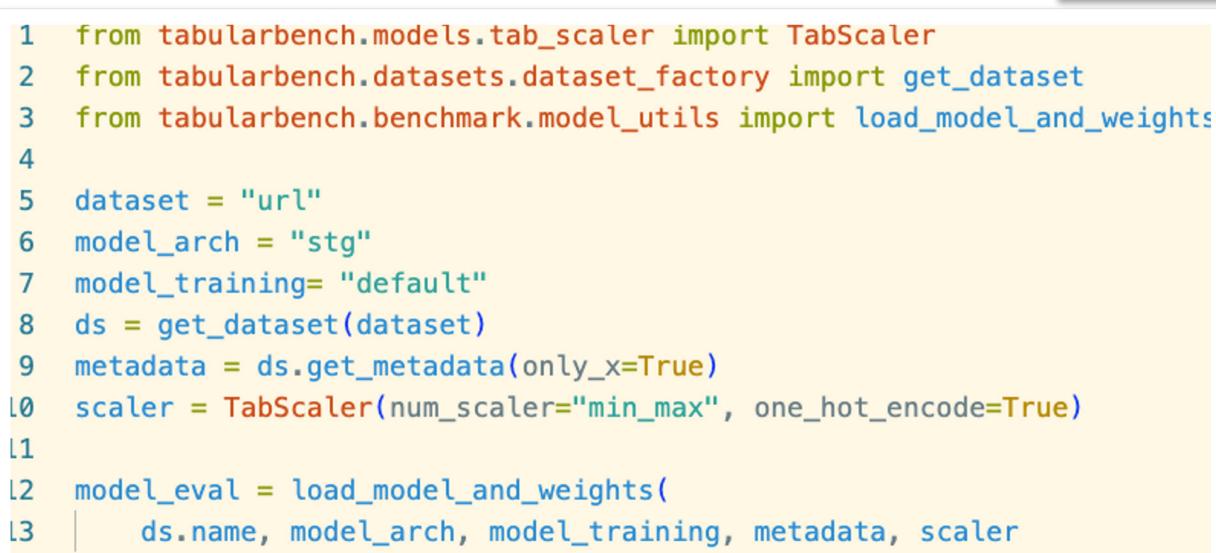
pip install tabularbench

### Models on Hugging Face



serval-uni-lu/tabularbench

# Access to datasets and models ↓



### **Constraints access and definition** ↓

2 from tabularbench.datasets.samples.lcld import get\_relation\_constraints 4 lcld\_constraints = get\_relation\_constraints() 6 new\_constraint = ( (Feature("term") == Constant(36)) (Feature("term") == Constant(48)) (Feature("term") == Constant(60))

# Learn more



**Paper** 



290 Pre-trained models

11 lcld\_constraints[3] = new\_constraint